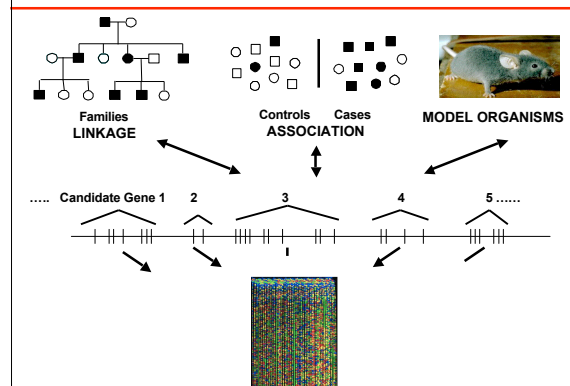


Medical Resequencing

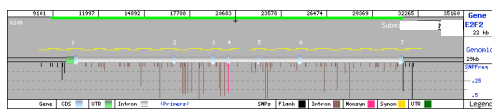
Debbie Nickerson

Department of Genome Sciences
University of Washington

Genetic Studies

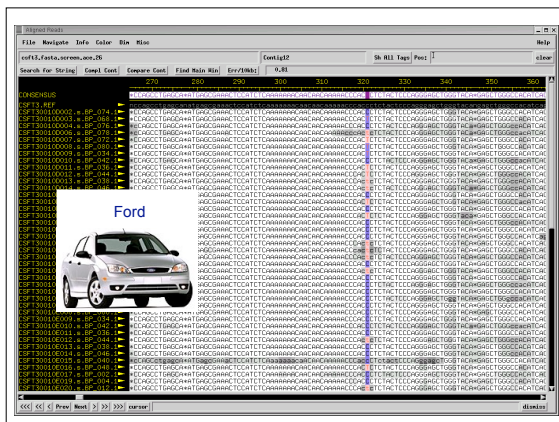
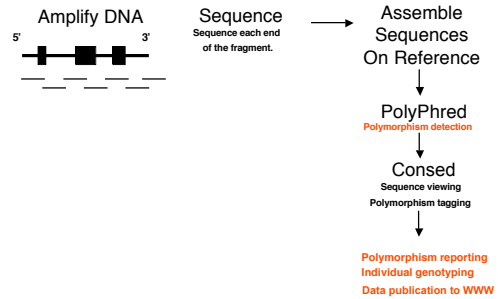


Overview of a Candidate Gene



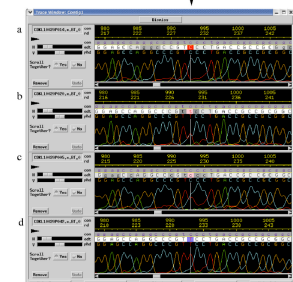
Average Gene Size - 26.5 kb ~ Compare 2 haploid - 1 in 1,200 bp
~130 SNPs (200 bp) - 15,000,000 SNPs
~ 44 SNPs \geq 0.05 MAF (600 bp) - 6,000,000 SNPs

Sequencing production and data analysis pipeline



Aston-Martin of SNP Detection - PolyPhred 5.0

* Matthew Stephens
Peggy Dyer-Robertson
Jim Sloan



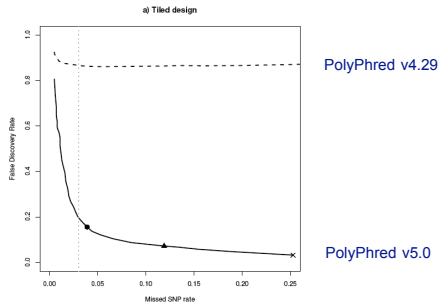
Comparison PolyPhred v4.29 versus v5.0

a) Tiled design

Missed SNP rate	PolyPhred v4.29 FDR	PolyPhred v5.0 FDR
0.00	~0.95	~0.95
0.01	~0.90	~0.60
0.02	~0.88	~0.25
0.03	~0.85	~0.15
0.05	~0.85	~0.12
0.10	~0.85	~0.08
0.15	~0.85	~0.06
0.20	~0.85	~0.05
0.25	~0.85	~0.05

PolyPhred v4.29

PolyPhred v5.0



PolyPhred 5 Scores - Provide Quantitative Assessment of SNP Genotype

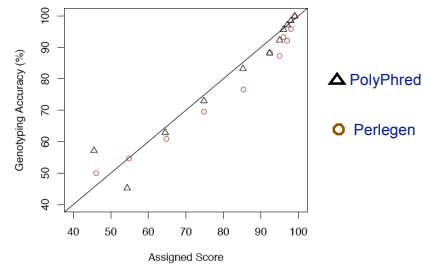
A scatter plot comparing the genotyping accuracy of PolyPhred (triangles) and Perlegen (circles) for 100 SNPs. The x-axis is labeled 'Assigned Score' and the y-axis is labeled 'Genotyping Accuracy (%)', both ranging from 40 to 100. A diagonal line represents perfect agreement (y=x). PolyPhred data points are clustered tightly along this line, indicating high accuracy. Perlegen data points are more scattered, with several points showing lower accuracy than PolyPhred's.

Assigned Score	PolyPhred Accuracy (%)	Perlegen Accuracy (%)
45	58	50
55	45	55
65	62	60
75	72	68
85	82	75
95	92	88
98	95	90
100	100	98

Legend:

- △ PolyPhred
- Perlegen

Double-Coverage - Automation = 93% of all SNPs, 100% of high-frequency SNPs, with no false positive SNPs identified, and 99.9% genotyping accuracy.



Double-Coverage - Automation = 93% of all SNPs, 100% of high-frequency SNPs, with no false positive SNPs identified, and 99.9% genotyping accuracy.

Comparison PolyPhred v5.0 to others

Final Discovery Rate

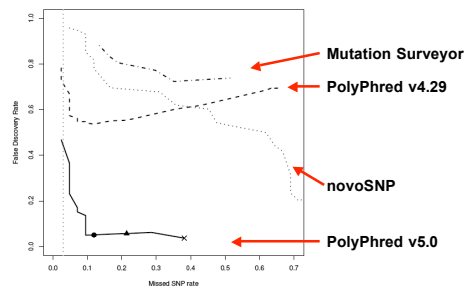
Missed SNP rate

Mutation Surveyor

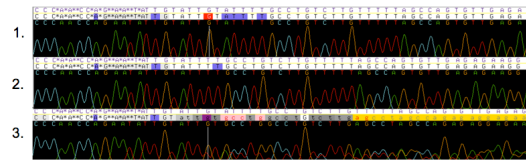
PolyPhred v4.29

novoSNP

PolyPhred v5.0

[illegible]

95% less than 15 bp



Bhangale et al (2005) Hum Mol Genet. 14: 59-69

Importance of short indels

- Indels are common and in LD with substitutions and can be used to improve the marker densities
- Indels are overrepresented as disease-causing mutations
 - ~24% of mutations in the HGMD are indels

type	no. of entries	percent
missense	1191	32.74%
deletion/insertion	4011	51.54%
frameshift	435	1.03%
regulatory	250	0.63%
stop	250	0.63%
small insertions	388	0.94%
small indels	388	0.94%
missense	86	0.20%
frameshift	388	0.94%
stop	388	0.94%
complex rearr / inversion	491	1.16%
gross deletions	235	0.59%
gross deletions	235	0.59%
total	3616	100.00%

24.22%

- Indels are common and in LD with substitutions and can be used to improve the marker densities
- Indels are overrepresented as disease-causing mutations
 - ~24% of mutations in the HGMD are indels

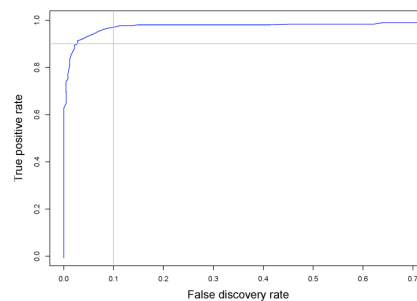
[illegible]

Indel-Detection Accuracy

The ROC curve illustrates the trade-off between the True Positive Rate (Y-axis) and the False Discovery Rate (X-axis). The curve starts at (0,0) and rises sharply, reaching a True Positive Rate of approximately 0.95 at a False Discovery Rate of 0.05. Beyond this point, the curve plateaus, maintaining a True Positive Rate near 1.0 for False Discovery Rates up to 0.7. A vertical grey line is drawn at a False Discovery Rate of 0.1, and a horizontal grey line is drawn at a True Positive Rate of approximately 0.9.

False discovery rate	True positive rate
0.0	0.0
0.01	0.62
0.02	0.75
0.03	0.82
0.04	0.88
0.05	0.92
0.06	0.94
0.07	0.95
0.10	0.96
0.20	0.97
0.30	0.97
0.40	0.97
0.50	0.97
0.60	0.98
0.70	0.98

For Every 9 True Positives - 1 False- Positives

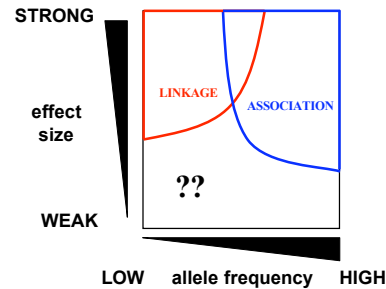


For Every 9 True Positives - 1 False- Positives

Medical Resequencing

- Discovery of rare functional variants -
 - Sequencing at the tails of the distribution
- Testing the Common Disease Common Variant (CDCV) hypothesis
 - Candidate genes very feasible
- Whole Genome Sequencing

Genetic Strategy Determined by Effect Size & Allele Frequency



Ardlie, Kruglyak & Seielstad (2002) Nat. Genet. Rev. 3: 299-309
Zondervan & Cardon (2004) Nat. Genet. Rev. 5: 89-100

ABCA1 and HDL-C

	Sequence variants unique to one group			
	Low HDL-C		High HDL-C	
	NS	S	NS	S
ABCA1	14	6	2	5
APOA1	1	0	0	1
LCAT	0	1	1	0
Canadians				
ABCA1	14	2	2	3
APOA1	0	1	0	0
LCAT	6	1	0	0

—Cohen et al, Science
305, 869-872, 2004

- Observed excess of rare, nonsynonymous variants in low HDL-C samples at ABCA1
- Demonstrated functional relevance in cell culture

Rare coding variants

- No single variant frequent enough for significant association
- Indications of function
 - Ratio of synonymous to nonsynonymous
 - Predicted function from evolutionary data
 - Wet bench tests

Medical Resequencing

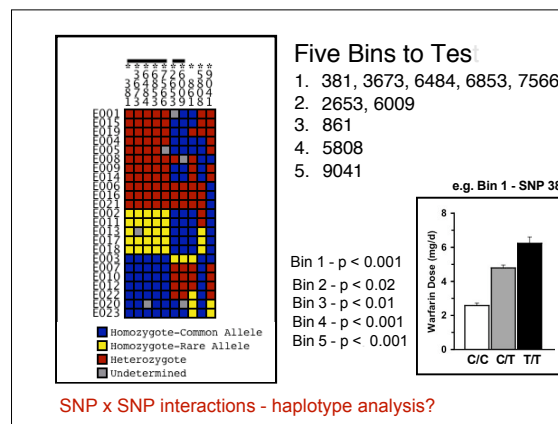
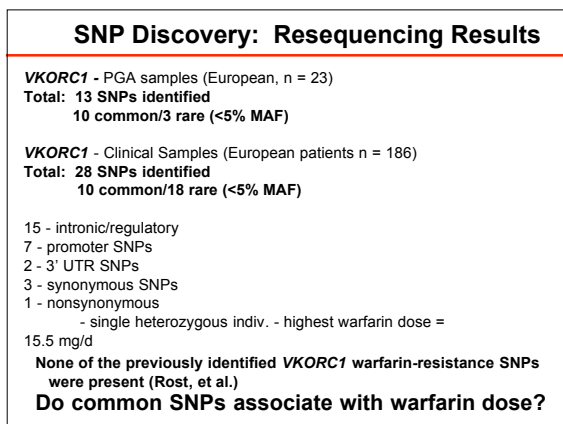
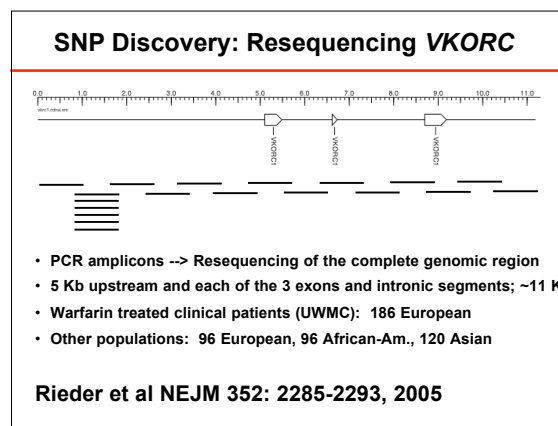
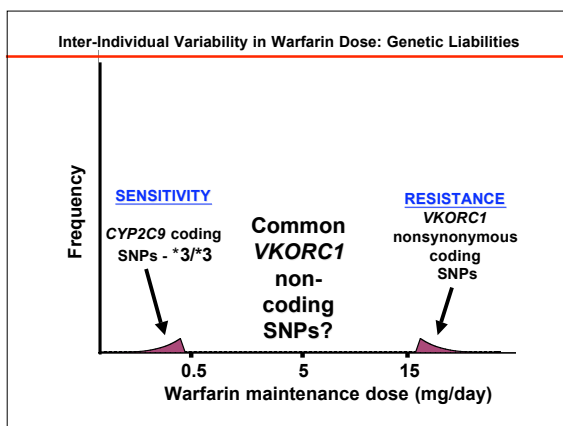
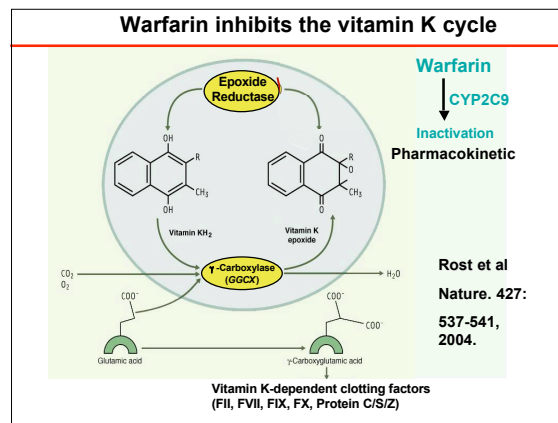
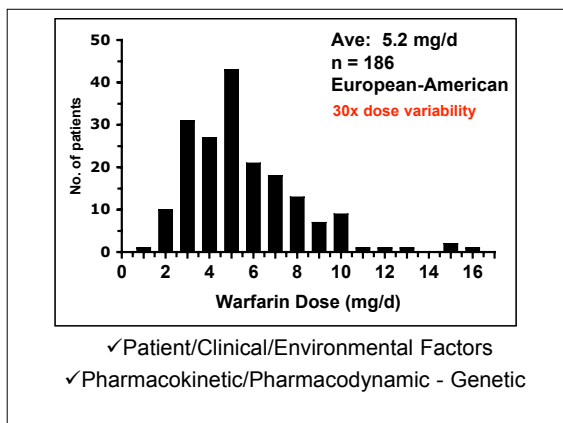
- Testing the Common Disease Common Variant (CDCV) hypothesis
 - Candidate genes very feasible
- What about rare variants (CDRV)?
- Whole genome using tagSNPs feasible but sequencing could be in the future

Warfarin Background

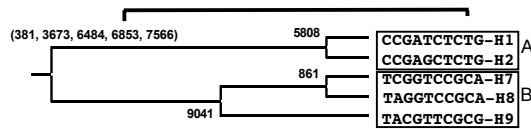
- Commonly prescribed oral anti-coagulant and acts as an inhibitor of the vitamin K cycle
- In 2003, 21.2 million prescriptions were written for warfarin (Coumadin™)
- Prescribed following MI, atrial fibrillation, stroke, venous thrombosis, prosthetic heart valve replacement, and following major surgery
- Difficult to determine effective dosage
 - Narrow therapeutic range
 - Large inter-individual variation



WARF+coumarin



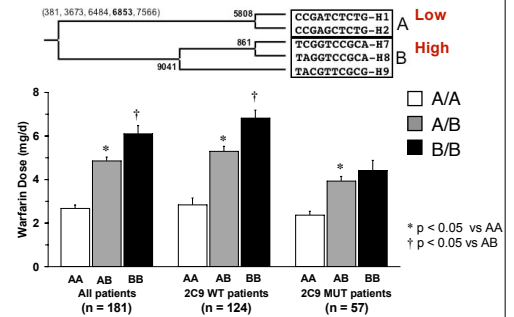
VKORC1 haplotypes cluster into divergent clades



Patients were assigned a clade diplotype:

e.g. Patient 1 - H1/H2 = A/A
Patient 2 - H1/H7 = A/B
Patient 3 - H7/H9 = B/B

VKORC1 clade diplotypes show a strong association with warfarin dose



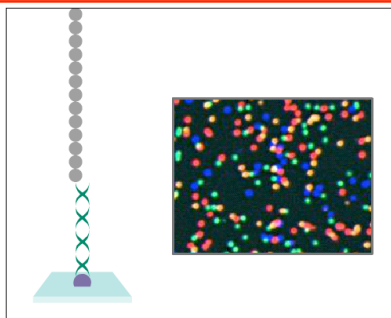
Medical Resequencing

- Discovery of rare functional variants -
 - Sequencing at the tails of the distribution
- Testing the Common Disease Common Variant (CDCV) hypothesis
 - Candidate genes very feasible
- Whole Genome Sequencing

SNP Genotyping -

Is it an intermediate stop on the way to whole-genome sequencing?

Long term sequencing - In situ approaches



Solexa - an example

Sequencing could be the ultimate genotyping tool

- More applications
- Further Technology Development